

Tema 60

Parámetros estadísticos. Cálculo, significado y propiedades

60.1 Parámetros estadísticos de una variable estadística

60.1.1 Introducción

Los dos diagramas “diferencial e integral” definidos en el tema 59 nos dan una representación visual de las variables estadísticas. La vista saca sucesivamente dos impresiones:

- por lectura de la escala de abscisas, retiene los valores de la variable situados en el centro de la distribución: valores de tendencia central.
- por lectura de las desviaciones a la tendencia central, una mayor o menor fluctuación alrededor de los valores centrales: es lo que se llama dispersión.

Estos dos elementos sintéticos de una distribución estadística permiten comparaciones.

Los resúmenes cuantitativos que informan de la variable los llamaremos estadísticos o características estadísticas.

60.1.2 Características de tendencia central

Media aritmética

Definición 1 *La media aritmética de una variable se define como la suma de los valores de la variable multiplicados por su frecuencia relativa correspondiente. Si llamamos X a la variable y x_i a sus valores, la media aritmética la notaremos*

por \bar{x} y se define por:

$$\bar{x} = \sum_{i=1}^k f_i x_i = \frac{1}{n} \sum_{i=1}^k n_i x_i$$

Si la variable estadística es discreta se aplica directamente la fórmula, pero cuando es continua se toman como valores x_i las marcas de clase de los intervalos y como frecuencias las de los intervalos.

Propiedades (1) La media de las diferencias a la media es cero, es decir,

$$\sum_{i=1}^k f_i (x_i - \bar{x}) = 0$$

Demostración:

$$\sum_{i=1}^k f_i (x_i - \bar{x}) = \sum_{i=1}^k f_i x_i - \sum_{i=1}^k f_i \bar{x} = \bar{x} - \bar{x} \sum_{i=1}^k f_i = \bar{x} - \bar{x} = 0 \quad C.Q.D. \square$$

(2) La media de los cuadrados de las desviaciones respecto de la media es mínima.

Demostración:

Sea $Q(a) = \sum_{i=1}^k f_i (x_i - a)^2$. Veamos para que valor de a es mínima $Q(a)$.

Para ello derivamos e igualamos a cero:

$$\begin{aligned} \frac{d}{da} Q(a) &= -2 \sum_{i=1}^k f_i (x_i - a) = 0 \Leftrightarrow \sum_{i=1}^k f_i (x_i - a) = 0 \Leftrightarrow \\ \Leftrightarrow 0 &= \sum_{i=1}^k f_i x_i - a \sum_{i=1}^k f_i = \bar{x} - a \Leftrightarrow a = \bar{x} \quad C.Q.D. \square \end{aligned}$$

A la cantidad $\sum_{i=1}^k f_i (x_i - \bar{x})^2$ se le llama varianza de la variable X .

(3) Si multiplicamos todos los valores x_i de una variable X por una cte k , entonces la media queda multiplicada por k .

Demostración:

Tenemos que probar que $\overline{kx} = k\bar{x}$.

$$\overline{kx} = \sum_{i=1}^k f_i (kx_i) = k \left(\sum_{i=1}^k f_i x_i \right) = k\bar{x} \quad C.Q.D. \square$$

Media geométrica, cuadrática y armónica

Definición 2 Se define la media geométrica de la variable estadística X por:

$$G = \sqrt[n]{x_1^{n_1} \dots x_k^{n_k}} = \sqrt[n]{\prod_{i=1}^k x_i^{n_i}}$$

Utilizando logaritmos y antilogaritmos se puede dar una expresión más operativa para calcular la media geométrica:

$$\begin{aligned} \log G &= \log \left(\prod_{i=1}^k x_i^{n_i} \right)^{\frac{1}{n}} = \frac{1}{n} \log \left(\prod_{i=1}^k x_i^{n_i} \right) = \frac{1}{n} \sum_{i=1}^k \log(x_i^{n_i}) = \frac{1}{n} \sum_{i=1}^k n_i \log x_i \Rightarrow \\ &\Rightarrow G = \sqrt[n]{\prod_{i=1}^k x_i^{n_i}} = \text{antilog} \left(\sum_{i=1}^k \frac{n_i \ln x_i}{n} \right) \end{aligned}$$

donde \log es la función logaritmo natural y antilog es la función exponencial natural (de base e).

Definición 3 Se define la media cuadrática de la variable estadística X por:

$$Q = \sqrt{\sum_{i=1}^k f_i x_i^2}$$

Definición 4 Se define la media armónica de la variable estadística X por:

$$H = \frac{1}{\sum_{i=1}^k f_i \frac{1}{x_i}}$$

La **relación** existente entre las medias es la siguiente (como se puede comprobar fácilmente):

$$H < G < \bar{x} < Q$$

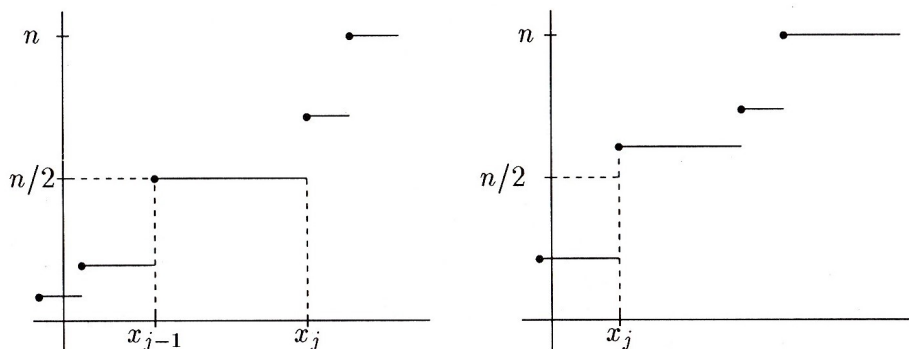
Mediana

Definición 5 Se define la mediana como aquel valor de la variable tal que, supuestos ordenados los valores de ésta en orden creciente, la mitad son menores o iguales y la otra mitad son mayores o iguales que dicho valor.

Cálculo Datos sin agrupar:

La figura siguiente correspondiente a un diagrama de frecuencias absolutas

acumuladas, recoge las dos situaciones que se pueden presentar.



Si la situación es como la de la figura de la derecha, es decir, si

$$N_{j-1} < \frac{n}{2} < N_j$$

entonces la mediana es

$$M_e = x_j$$

Si la situación que se presenta es como la de la figura de la izquierda, entonces la mediana queda indeterminada, aunque en este caso se toma como mediana la media aritmética de los dos valores entre los que se produce la indeterminación; así pues, si

$$N_{j-1} = \frac{n}{2} < N_j$$

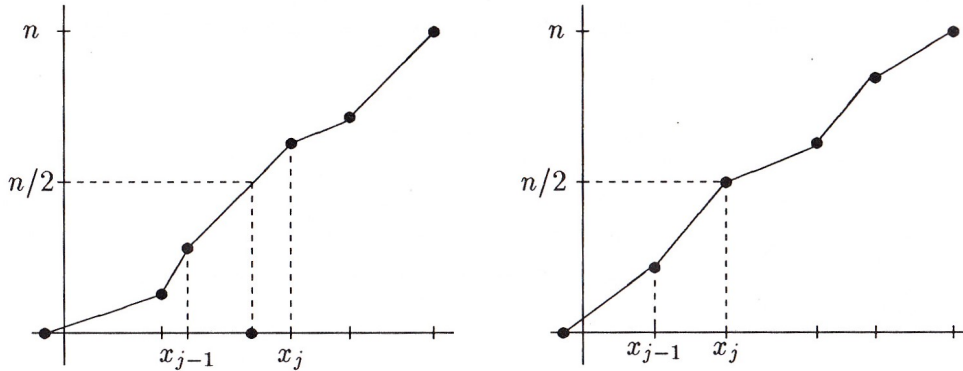
entonces la mediana es

$$M_e = \frac{x_{j-1} + x_j}{2}$$

Datos agrupados:

La figura siguiente, correspondiente a un polígono de frecuencias absolutas

acumuladas, nos plantea de nuevo dos situaciones diferentes a considerar



El más sencillo, el de la derecha, en el que existe una frecuencia absoluta acumulada N_j tal que $\frac{n}{2} = N_j$, la mediana es:

$$M_e = x_j$$

Si la situación es como la que se presenta en la figura de la izquierda, en la que

$$N_{j-1} < \frac{n}{2} < N_j$$

entonces, la mediana, está en el intervalo $[x_{j-1}, x_j[$, es decir, entre x_{j-1} y x_j , tomándose en ese caso, por razonamientos de proporcionalidad, como mediana el valor

$$M_e = x_{j-1} + \frac{\frac{n}{2} - N_{j-1}}{n_j} c_j$$

siendo $c_j = x_j - x_{j-1}$.

Propiedades 1) La mediana sólo depende del orden en que están colocados los valores y no de los valores en sí.

2) La mediana de una variable discreta es siempre uno de los valores posibles de la variable y es más concreta que la media.

3) Es respecto de la mediana cuando la desviación absoluta media (concepto que se definirá más adelante) es mínima.

Moda

Definición 6 La moda se define como aquel valor de la variable al que corresponde máxima frecuencia (absoluta o relativa).

La moda es fácil de calcular y posee un significado concreto excelente. Además, algunas distribuciones pueden presentar varias modas (son las denominadas plurimodales), que corresponden cada una a un máximo local del diagrama diferencial.

Para calcularla también será necesario distinguir si los datos están agrupados o no.

Cálculo Datos sin agrupar:

Para datos sin agrupar, la determinación del valor o valores (ya que puede haber más de uno) modal/es es muy sencillo. Basta observar a que valor le corresponde una mayor n_i . Este será la moda.

Datos agrupados:

Si los datos se presentan agrupados en intervalos es necesario, a su vez, distinguir si éstos tienen o no igual amplitud.

Si tienen amplitud cte c , una vez identificado el intervalo modal $[x_{j-1}, x_j]$, es decir, el intervalo al que corresponde mayor frecuencia absoluta

$$n_j = \max(n_1, \dots, n_k)$$

la moda se define, también por razones geométricas como

$$M_o = x_{j-1} + \frac{n_{j+1}}{n_{j-1} + n_{j+1}}c$$

Si los intervalos tuvieran distinta amplitud c_j , primero debemos normalizar las frecuencias absolutas n_j , determinando los cocientes

$$l_j = \frac{n_j}{c_j} \quad \forall j = 1, \dots, k$$

y luego aplicar la regla definida para el caso de intervalos de amplitud constante a los l_j . Es decir, primero calcular el $l_j = \max(l_1, \dots, l_k)$ para determinar el intervalo modal $[x_{j-1}, x_j]$ y luego aplicar la fórmula

$$M_o = x_{j-1} + \frac{l_{j+1}}{l_{j-1} + l_{j+1}}c_j$$

siendo $c_j = x_j - x_{j-1}$.

Propiedad: Relación entre la media, la moda y la mediana

Generalmente, la mediana está comprendida entre la moda y la media, es decir,

$$M_o \leq M_e \leq \bar{x}$$

60.2 Medidas de posición no central

60.2.1 Cuantiles

Definición 7 *El cuantil $p_{r/k}$, $r = 1, \dots, k - 1$ se define como aquel valor de la variable que divide la distribución de frecuencias, previamente ordenada de forma creciente, en dos partes, estando el $(100 \frac{r}{k})\%$ de ésta formado por valores menores que $p_{r/k}$.*

Si $k = 4$ los (tres) cuantiles reciben el nombre cuartiles. Si $k = 10$ los (nueve) cuantiles reciben, en este caso, el nombre de deciles. Por último, si $k = 100$ los (noventa y nueve) cuantiles reciben el nombre de centiles.

Obsérvese que siempre que r y k mantengan la misma proporción $\left(\frac{r}{k}\right)$ obtendremos el mismo valor. Es decir, por ejemplo, el primer cuartil es igual al vigésimo quinto centil.

En este sentido, la mediana M_e es el segundo cuartil, o el quinto decil, etc.

Para el cálculo de los cuantiles de nuevo hay que considerar si los datos vienen o no agrupados en intervalos.

Datos sin agrupar:

Si los datos vienen sin agrupar y

$$N_{j-1} < \frac{r}{k}n < N_j$$

entonces, el r -ésimo cuantil de orden k será

$$p_{r/k} = x_j$$

valor al que corresponde la frecuencia absoluta acumulada N_j .

Si la situación fuera de la forma

$$N_{j-1} = \frac{r}{k}n < N_j$$

tomaríamos, en esta situación indeterminada,

$$p_{r/k} = \frac{x_{j-1} + x_j}{2}$$

Datos agrupados:

Si los datos se presentan agrupados y, para algún j , fuera

$$\frac{r}{k}n = N_j$$

entonces, el r -ésimo cuantil de orden k sería

$$p_{r/k} = x_j$$

Por último, si fuera

$$N_{j-1} < \frac{r}{k}n < N_j$$

el intervalo a considerar sería $[x_{j-1}, x_j[$, al que corresponde frecuencia absoluta n_j y absoluta acumulada N_j , siendo entonces el cuantil el dado por la expresión,

$$p_{r/k} = x_{j-1} + \frac{\frac{r}{k}n - N_{j-1}}{n_j}c_j \quad \forall r = 1, \dots, k - 1$$

en donde c_j es la amplitud del intervalo $[x_{j-1}, x_j[$.

Si el intervalo a considerar fuera el $[x_0, x_1[$, se tomaría en la expresión anterior $N_{j-1} = 0$.

60.3 Características de dispersión

60.3.1 Introducción

Las medidas de posición estudiadas en las secciones anteriores servían para resumir la distribución de frecuencias de un sólo valor. Las medidas de dispersión, a las cuales dedicaremos esta sección, tienen como propósito estudiar lo concentrada que está la distribución en torno a algún promedio.

60.3.2 Características de dispersión absolutas

Varianza

Definición 8 *Es la media aritmética de los cuadrados de las desviaciones a la media, es decir,*

$$\sigma_X^2 \equiv V(X) \equiv \text{Var}(X) = \sum_{i=1}^k f_i (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^k n_i (x_i - \bar{x})^2$$

(donde si la variable es continua, x_i representa las marcas de clase).

Veamos algunas **propiedades** de la varianza:

(1) $\sigma_X^2 \geq 0$

La varianza tomará el valor cero si la distribución presenta un único valor, en cuyo caso $\bar{x} = x_1$, y por tanto la distribución no presenta dispersión.

(2) La varianza es la medida de dispersión óptima

(3) Si a los valores de una variable les sumamos una cte, la varianza de la nueva variable así construida coincide con la varianza de la variable original.

Demostración:

Consideramos la variable X con valores x_i y construimos la variable Y con valores $y_i = x_i + b \quad \forall i = 1, \dots, k$, y $b \in \mathbb{R}$. Tenemos que demostrar que $\sigma_X^2 = \sigma_Y^2$.

$$\sigma_Y^2 = \sum_{i=1}^k f_i (y_i - \bar{y})^2 = \sum_{i=1}^k f_i (x_i + b - \bar{x} + b)^2 = \sum_{i=1}^k f_i (x_i - \bar{x})^2 = \sigma_X^2 \quad \square$$

(4) Si todos los valores de una variable X son multiplicados por $a \in \mathbb{R}$, entonces $\sigma_Y^2 = a^2 \sigma_X^2$ donde $Y = aX$.

Demostración:

Consideramos la variable X con valores x_i y construimos la variable Y con valores $y_i = ax_i + b \quad \forall i = 1, \dots, k$, y $b \in \mathbb{R}$. Tenemos que demostrar que $\sigma_Y^2 = a^2 \sigma_X^2$.

$$\begin{aligned} \sigma_Y^2 &= \sum_{i=1}^k f_i (y_i - \bar{y})^2 = \sum_{i=1}^k f_i (ax_i - a\bar{x})^2 = \sum_{i=1}^k f_i [a(x_i - \bar{x})]^2 = \\ &= \sum_{i=1}^k f_i a^2 (x_i - \bar{x})^2 = a^2 \sum_{i=1}^k f_i (x_i - \bar{x})^2 = a^2 \sigma_X^2 \quad \square \end{aligned}$$

(5) **Teorema de König**

$$\sigma_X^2 = \sum_{i=1}^k f_i x_i^2 - \bar{x}^2$$

Demostración:

$$\begin{aligned} \sigma_X^2 &= \frac{1}{n} \sum_{i=1}^k n_i (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^k n_i (x_i^2 - 2x_i \bar{x} + \bar{x}^2) = \\ &= \frac{1}{n} \sum_{i=1}^k n_i x_i^2 - 2\bar{x} \sum_{i=1}^k \frac{n_i}{n} x_i + \frac{1}{n} \sum_{i=1}^k n_i \bar{x}^2 = \\ &= \frac{1}{n} \sum_{i=1}^k n_i x_i^2 - 2\bar{x}^2 + \bar{x}^2 = \sum_{i=1}^k f_i x_i^2 - \bar{x}^2 \quad C.Q.D. \square \end{aligned}$$

Desviación típica

Definición 9 Es la raíz cuadrada positiva de la varianza, es decir,

$$\sigma_X = +\sqrt{\sigma_X^2} = +\sqrt{\sum_{i=1}^k f_i (x_i - \bar{x})^2}$$

La desviación típica no tiene un sentido muy concreto en sí misma. Tiene significado sólo para comparar distribuciones.

60.3.3 Características de dispersión relativas

Vienen dadas en forma de coeficientes adimensionales, lo que va a favorecer la comparación entre los grados de dispersión de distribuciones heterogéneas.

Coefficiente de variación de Pearson

Definición 10 Es el coeficiente

$$CV = \frac{\sigma_X}{\bar{x}}$$

Es una cantidad sin dimensión, independiente de las unidades elegidas, es decir, invariante si se efectúa un cambio de escala. Por tanto, permite comparar distribuciones., y en este sentido, cuanto más pequeño sea dicho coeficiente, más homogénea será una variable al compararla con otra.

Otra propiedad muy importante de este coeficiente es la siguiente:

- Si $\bar{x} < \sigma \Rightarrow \bar{x}$ no tiene representatividad ninguna
- Si $\sigma = 0 \Rightarrow CV = 0 \Rightarrow$ hay máxima representatividad de la media (no hay dispersión)

60.4 MOMENTOS

60.4.1 Momentos no centrales y centrales

Definición 11 Sea $r \in \mathbb{N} \cup \{0\}$ y $a \in \mathbb{R}$. Se llama momento de orden “ r ” respecto de “ a ” a la cantidad

$${}_a m_r = \sum_{i=1}^k f_i (x_i - a)^r = \frac{1}{n} \sum_{i=1}^k n_i (x_i - a)^r$$

Según los valores de a se definen varias clases de momentos:

Definición 12 Se llama momento no central de orden r ó respecto de $a = 0$ a la cantidad

$$m_r = \sum_{i=1}^k f_i x_i^r$$

Definición 13 Se llama momento central de orden r ó respecto de $a = \bar{x}$ a la cantidad

$$\mu_r = \sum_{i=1}^k f_i (x_i - \bar{x})^r$$

Los primeros momentos de ambos tipos son:

$$\begin{aligned} m_0 &= \sum_{i=1}^k f_i = 1 & \mu_0 &= 1 \\ m_1 &= \sum_{i=1}^k f_i x_i = \bar{x} & \mu_1 &= \sum_{i=1}^k f_i (x_i - \bar{x}) = 0 \\ m_2 &= \sum_{i=1}^k f_i x_i^2 = \sigma^2 + \bar{x}^2 & \mu_2 &= \sum_{i=1}^k f_i (x_i - \bar{x})^2 = \sigma^2 \end{aligned}$$

60.4.2 Relación entre los momentos

Proposición 14 Momentos centrales en función de los no centrales

$$\mu_r = \sum_{t=0}^{r-2} (-1)^t \binom{r}{t} m_1^t m_{r-t} + (-1)^{r-1} m_1^r (r-1)$$

Demostración:

$$\begin{aligned}
\mu_r &= \sum_{i=1}^k f_i \left(x_i - \underbrace{m_1}_{=\bar{x}} \right)^r = \sum_{i=1}^k f_i \sum_{t=0}^r (-1)^t \binom{r}{t} m_1^t x_i^{r-t} = \\
&= \sum_{t=0}^r (-1)^t \binom{r}{t} m_1^t \underbrace{\sum_{i=1}^k f_i x_i^{r-t}}_{=m_{r-t}} = \\
&= \sum_{t=0}^{r-2} (-1)^t \binom{r}{t} m_1^t m_{r-t} + (-1)^{r-1} \underbrace{\binom{r}{r-1}}_{=r} \underbrace{m_1^{r-1} m_1}_{=m_1^r} + (-1)^r \underbrace{\binom{r}{r}}_{=1} m_1^r \underbrace{m_0}_{=1} = \\
&= \sum_{t=0}^{r-2} (-1)^t \binom{r}{t} m_1^t m_{r-t} + (-1)^{r-1} r m_1^r + (-1)^r m_1^r = \\
&= \sum_{t=0}^{r-2} (-1)^t \binom{r}{t} m_1^t m_{r-t} + (-1)^{r-1} m_1^r (r-1) \quad C.Q.D. \square
\end{aligned}$$

Proposición 15 *Momentos no centrales en función de los centrales*

$$m_r = \sum_{t=0}^{r-2} \binom{r}{t} m_1^t \mu_{r-t} + m_1^r$$

Demostración:

$$\begin{aligned}
m_r &= \sum_{i=1}^k f_i x_i^r = \sum_{i=1}^k f_i [(x_i - m_1) + m_1]^r = \sum_{i=1}^k f_i \sum_{t=0}^r \binom{r}{t} m_1^t (x_i - m_1)^{r-t} = \\
&= \sum_{t=0}^r \binom{r}{t} m_1^t \underbrace{\sum_{i=1}^k f_i (x_i - m_1)^{r-t}}_{=\mu_{r-t}} = \\
&= \sum_{t=0}^{r-2} \binom{r}{t} m_1^t \mu_{r-t} + \binom{r}{r-1} m_1^{r-1} \underbrace{\mu_1}_{=0} + \binom{r}{r} m_1^r \underbrace{\mu_0}_{=1} = \\
&= \sum_{t=0}^{r-2} \binom{r}{t} m_1^t \mu_{r-t} + m_1^r \quad C.Q.D. \square
\end{aligned}$$

Así, para $r = 2, 3, 4$ tenemos:

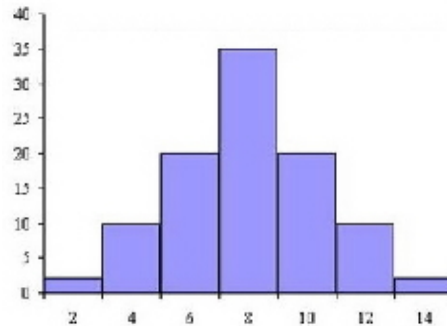
$$\begin{array}{ll}
\mu_2 = m_2 - m_1^2 \text{ (Teorema de König)} & m_2 = \mu_2 + m_1^2 \text{ (Teorema de König)} \\
\mu_3 = m_3 - 3m_2m_1 + 2m_1^3 & m_3 = \mu_3 + 3\mu_2m_1 + m_1^2 \\
\mu_4 = m_4 - 4m_3m_1 + 6m_1^2m_2 - 3m_1^4 & m_4 = \mu_4 + 4\mu_3m_1 + 6\mu_2m_1^2 + m_1^2
\end{array}$$

60.5 CARACTERÍSTICAS DE FORMA

Además de la tendencia central y de la dispersión, se puede tratar de caracterizar la forma de una distribución mediante índices resumidos: índices de asimetría y de aplastamiento. Estos índices nos indican cual es la forma de distribución de la variable.

La simetría estudia la falta de simetría de la variable X respecto del eje vertical $x = \bar{x}$.

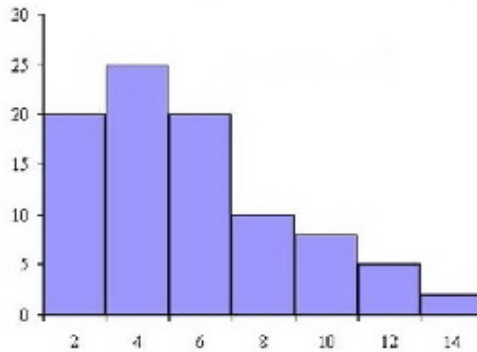
Definición 16 Una distribución de frecuencias se dice simétrica si lo es su representación gráfica o, lo que es lo mismo, cuando son iguales las frecuencias correspondientes a valores de la variable equidistantes de una valor central.



Si una distribución es simétrica, se verifica tanto en el caso discreto como en el continuo que:

$$\bar{x} = M_o = M_e$$

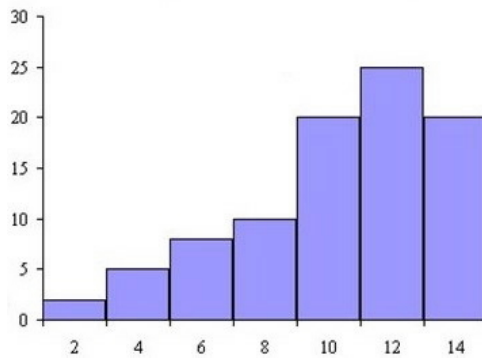
Definición 17 Diremos que una distribución es asimétrica a la derecha o positiva cuando las frecuencias descienden más lentamente por la derecha que por la izquierda o, lo que es lo mismo, en la representación gráfica la rama de la derecha es más larga que la de la izquierda respecto de la moda.



En este caso se verifica que:

$$\bar{x} > M_o$$

Definición 18 Diremos que una distribución es asimétrica a la izquierda o negativa cuando las frecuencias descienden más lentamente por la izquierda que por la derecha o, lo que es lo mismo, en la representación gráfica la rama de la izquierda es más larga que la de la derecha respecto de la moda.



En este caso se verifica que:

$$\bar{x} < M_o$$

60.5.1 Coeficientes de asimetría¹

Coefficiente de FISHER

$$\gamma_1(X) = \sum_{i=1}^k f_i \left(\frac{x_i - \bar{x}}{\sigma} \right)^3 = \frac{\mu_3}{\sigma^3} = \frac{\mu_3}{\mu_2^{\frac{3}{2}}}$$

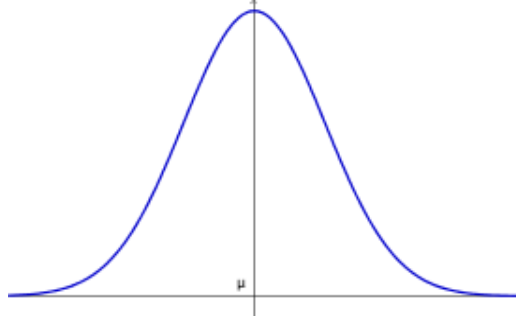
Este coeficiente no tiene dimensión, es invariante por cambio de origen y de escala, y nulo para las distribuciones simétricas.

60.5.2 Coeficientes de aplastamiento

Se define la curtosis o aplastamiento de una variable mediante la comparación de dicha variable con la variable patrón, denominada variable normal o de Laplace-Gauss, y mide la agrupación de los valores de la variable en torno a la media

¹Si $\psi(X)$ es un coeficiente de asimetría de los que se van a definir a continuación, entonces se tiene el siguiente criterio de asimetría: si $\psi(X) > 0$ la asimetría es positiva; si $\psi(X) = 0$ es simétrica; y si $\psi(X) < 0$ la asimetría es negativa.

(cuanto más concentrados, mayor será el agrupamiento).

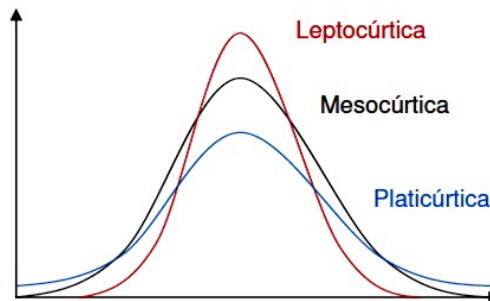


Coefficiente de curtosis de Pearson

$$\gamma_2(X) = \sum_{i=1}^k f_i \left(\frac{x_i - \bar{x}}{\sigma} \right)^4 - 3 = \frac{\mu_4}{\sigma^4} - 3 = \frac{\mu_4}{\mu_2^2} - 3$$

En las proximidades de la media si:

- (1) $\gamma_2(X) > 0 \Rightarrow$ la distribución es más apuntada que la normal y se denomina leptocúrtica.
- (2) $\gamma_2(X) = 0 \Rightarrow$ la distribución es igual de apuntada que la normal y se denomina mesocúrtica.
- (3) $\gamma_2(X) < 0 \Rightarrow$ la distribución es menos apuntada que la normal y se denomina platicúrtica.



Es un coeficiente sin dimensión e invariante por cambio de origen y de escala, en el que la cte 3 se elige de modo que el coeficiente sea nulo cuando la distribución es normal.

Debido a la desigualdad $\mu_4 \geq \sigma^4$, el coeficiente de aplastamiento es siempre mayor que -2, y es igual a -2 en el límite, es decir, cuando las observaciones x_i son iguales entre sí o cuando existen solamente dos valores posibles con frecuencias iguales.